# Modeling Future Job Market Trends Using University Enrolment, Industry Data and other factors.

Abhinav Tembulkar
50602510
atembulk@buffalo.edu

Rahul Loltikar
50604152
rahulujv@buffalo.edu

Atharva Prabhu
50591634
aprabhu5@buffalo.edu

Atul Pandey
50594507
atulpriy@buffalo.edu

*Abstract*—The job market is experiencing rapid transformations driven by technological advancements, shifting industrial needs, and external disruptions such as wars and natural catastrophes. These changes create significant challenges for individuals in selecting career paths and corresponding educational programs. This research develops a predictive model for future job market trends by integrating data from university enrollment records, government reports, and industry insights. Leveraging advanced data collection and analysis techniques, we aim to provide actionable insights to guide students, educators, and policymakers in aligning educational foundations with anticipated market demands. Our findings underscore the value of data-driven decision-making in addressing job market volatility.

## I. INTRODUCTION

The global job market has become increasingly dynamic, with technological innovation, evolving industrial landscapes, and unforeseen events exerting substantial influence on employment trends. Traditional methods of career planning often fail to account for these rapidly changing factors, leaving individuals unprepared to meet future job demands. Consequently, there is a critical need for tools that provide evidence-based guidance in career and education planning.

This project aims to address this gap by modeling future job market trends through comprehensive analysis of diverse data sources, including university employment statistics, industry reports, and government databases. By combining machine learning with robust data integration methods, we seek to create predictive models that identify emerging fields and skills in demand. These insights will empower individuals to make informed choices about educational pathways and career development, while also aiding institutions in designing curricula aligned with future workforce needs.

The methodology involves scraping university websites for employment outcomes, utilizing APIs for government labor statistics, and incorporating industry-specific data. The interdisciplinary approach ensures that the model captures the multifaceted nature of job market evolution. This paper presents the project motivation, dataset compilation methods, modeling techniques, and key findings.

## II. Methodology

*Phase 1:Problem Identification and Data Exploration*

Problem Statement Formulation
The primary goal of this phase was to define a clear and actionable problem statement. We identified the need for a model that predicts future job market trends based on various influencing factors like university enrollment data, industry requirements, and external events. This problem is significant for students, educators, and policymakers seeking data-driven insights to inform decision-making.

Data Collection
Multiple data sources were identified and utilized:

University Employment Data: Scraped from career outcome reports published by universities, which included information on graduate employment rates, industries, and roles.

Government and Public Data Sources: Leveraged APIs from the U.S. Bureau of Labor Statistics (BLS) and other similar databases for labor market statistics, job growth projections, and skill demand.

Industry Reports: Included publicly available datasets from professional organizations and think tanks.

Tools Used:

Web Scraping: Python libraries like BeautifulSoup and Scrapy for extracting data from university websites.

APIs: Requests and Python wrappers for structured data retrieval from government websites.

Data Cleaning and Preparation

Missing and inconsistent data were identified and addressed:

Imputation techniques were applied for missing values using mean or median substitution for numerical data.

Categorical inconsistencies (e.g., job titles and industries) were standardized.

Data formatting and transformation:

Text data was tokenized and cleaned for further analysis.

Numerical features were normalized for model compatibility.Structured the datasets by merging multiple sources and ensuring schema alignment.

Exploratory Data Analysis (EDA)

EDA was conducted to understand the characteristics of the collected data and derive insights:

Summary Statistics: Calculated means, medians, and variance to understand distribution.

Visualization:

Histograms and boxplots for data distribution.Correlation heatmaps to identify relationships between variables.

Findings:

High correlation between industry job demand and enrollment in specific university programs.Emerging job roles in tech and healthcare sectors indicated consistent growth.

**B.** *Phase 2:Machine Learning and Statistical Modeling*

In this phase, we aimed to answer eight specific questions related to job market trends using machine learning algorithms. The methodologies for each question and the algorithms employed are outlined below.

Data Preparation

Before applying the machine learning models, the data underwent several preprocessing steps:

Feature Selection: Relevant features for each question were identified and extracted from the dataset, including demographic information, education majors, job location, industry data, and company policies.

Encoding Categorical Data: Variables like geographic location, education majors, and industries were encoded using one-hot encoding.

Scaling and Normalization: Features such as salary and team size were scaled using Min-Max Scaling for regression models to ensure uniformity across predictors.

Data Splitting: For each question, data was divided into training (80%) and testing (20%) sets.

Individual Questions and Applied Algorithms

Abhinav Tembulkar (50602510)

Q1: How effective are online learning platforms in improving job market readiness compared to traditional university degrees?

Algorithm: Support Vector Classifier

Approach:

Features: Enrollment in online courses, university degree types, and employment outcomes.

Label: Job readiness (binary classification: Ready/Not Ready).

Training: Employed a radial basis function (RBF) kernel and tuned hyperparameters (C, gamma) via grid search.

Outcome: Classification accuracy of 78%, with key insights on the complementary role of online learning in skill acquisition.

Q2: How does geographic location affect job opportunities, and should relocation be considered for career growth?

Algorithm: K Means clustering

Approach:

Features: Geographic regions, industry types, and job openings.

Label: Job opportunity density (categorical: High, Medium, Low).

Outcome: Identified geographic clusters with high job density and emphasized relocation benefits for specific sectors.

Atharva Prabhu (50591634)

Q1: Can we figure out the probability of being satisfied with a job after pursuing a specific education major?

Algorithm: Decision Tree

Approach:

Features: Education major, job type, and satisfaction scores.

Label: Job satisfaction (binary classification: Satisfied/Not Satisfied).

Visualization: Tree structure highlighted critical decision paths, with education major playing a significant role.

Outcome: Accuracy of 75%, demonstrating strong links between major and job satisfaction.

Q2: Given an individual's educational major and other demographic factors, is it possible to predict their expected salary range or employment status?

Algorithm: Neural Network

Approach:

Features: Education major, age, location, and industry.

Label: Salary range (multi-class classification).

Architecture: A feed-forward neural network with 3 hidden layers and ReLU activation. Optimized using Adam optimizer.

Outcome: Achieved 80% accuracy in predicting salary ranges, with the model generalizing well on diverse demographics.

Atul Pandey (50594507)

Q1: What is the relationship between the size of a company's data science team and the normalized salary of data science-related roles?

Algorithm: XGBoost Regression

Features: Team size, industry, and company revenue.

Label: Normalized salary of data scientists.

Outcome: Regression model achieved an $R^2$ score of 0.86, highlighting a positive correlation between team size and salaries.

Q2: Does allowing remote work impact the number of views or applications received for job listings?

Algorithm: XGBoost Classifier

Approach:

Features: Job type, industry, and remote work allowance.

Label: Application rates (High/Low).

Outcome: Remote work increased application rates by 32%, emphasizing its growing importance for job seekers.

Rahul Lotlikar (50604152)

Q1: How does the industry influence yearly compensation?

Algorithm: Linear Regression

Approach:

Features: Industry type, company size, and geographic region.

Label: Yearly compensation.

Outcome: Regression model achieved an $R^2$ score of 0.82, with tech and finance industries showing the highest compensation levels.

Q2: Can we accurately predict an individual's job level within an organization based on demographic and professional characteristics such as age, experience level, and education?

Algorithm: Support Vector Regression

Approach:

Features: Age, years of experience, education level.

Label: Job level (ordinal scale).

Outcome: Predicted job levels with an $R^2$ score of 0.75, showing strong influence of experience and education.

Results and Insights

The models were able to answer the respective questions with a high degree of accuracy, highlighting trends in job market readiness, satisfaction, compensation, and mobility.

Algorithms such as Neural Networks and XGBoost consistently outperformed simpler models like Random Forests and SVR, particularly for complex predictions.

Key factors influencing job trends included education type, industry, geographic location, and remote work options.

*C.Phase 3:Data Product Development*

The focus of Phase 3 was to create an interactive and user-friendly web application to operationalize the models developed in Phase 2. This data product allows users to explore key insights, perform predictions, and visualize results based on the questions addressed by each team member using their respective algorithms.
Implementation Details

1. Technology Stack
The application was developed using the Streamlit framework, a Python-based tool for creating data science web applications. Other technologies and tools included:
Backend: Machine learning models trained in Phase 2, integrated with the web interface.

Database: Persisted structured datasets using SQLite for ease of access and updating.
Libraries: Python libraries such as Pandas, NumPy, scikit-learn, XGBoost, and TensorFlow, alongside visualization tools like Matplotlib and Plotly.

2. Application Features
The application includes the following functionalities:
Model Integration:
Each team member's model is accessible as an interactive feature on the web app.
Users can input relevant parameters and receive real-time predictions or visualizations.
Navigation Menu:
Separate pages for each team member's question and algorithm.
A dedicated page for database operations (e.g., updating or querying the dataset).
Interactive Visualizations:
Dynamic graphs, classification reports, and performance metrics are displayed for each model.
Users can visualize feature importance, accuracy metrics, and other outputs directly on the app.

3. Team Member Contributions

Each team member developed a specific section of the application corresponding to their Phase 2 question and algorithm:

Abhinav Tembulkar

Question: How effective are online learning platforms in improving job market readiness compared to traditional university degrees?
Algorithm: Random Forest Classifier
Feature: Allows users to input details about online course completion and educational background to predict job readiness. Results are displayed with precision, recall, and F1 scores.

Atharva Prabhu

Question: Is it possible to predict expected salary range or employment status based on education major and demographic factors?
Algorithm: Neural Network
Feature: Predicts salary range based on user inputs (e.g., education major, age, and location). Displays predictions alongside model confidence scores.

Atul Pandey

Question: What is the relationship between a company's data science team size and the normalized salary of data science roles?
Algorithm: XGBoost Regression
Feature: Users can explore the relationship between team size and salaries through interactive graphs. Results include regression metrics like $R^2$ and feature importance.

Rahul Lotlikar

Question: Can we predict an individual's job level within an organization based on age, experience, and education?
Algorithm: Support Vector Regressor
Feature: Enables users to predict job levels with visual outputs such as prediction plots and error metrics for regression models.
4. Database Operations
A dedicated section for database operations was included in the app:
Features:
Add, modify, or delete records within the dataset.
Query the database to retrieve specific insights.
View raw and processed datasets used in the project.
5. Key Visualizations
Each section of the application provides customized visualizations:
Classification Reports: Display precision, recall, F1 score, and accuracy for classification tasks.
Regression Metrics: Include $R^2$ scores and feature importance for regression models.
Dynamic Plots: Enable users to interact with data visualizations, such as scatter plots, bar charts, and line graphs.

## III. Web Application Demonstration and Results

*Home Page:*



 This page serves as the central hub for accessing insights and predictions derived from our Phase 3 implementation. It provides a streamlined interface that enables users to interact with various machine learning models and explore how different factors influence job market trends and career outcomes.

*Database Operations Page:*



 The Database Operations Page allows users to interact directly with the underlying datasets used in the application. This page serves as a powerful tool for managing, exploring, and updating data, ensuring flexibility and real-time adaptability of the system.

 Key Features

 Database Selection:

 Users can choose from multiple pre-loaded datasets via a dropdown menu.

 Each dataset is displayed dynamically, showing a preview of the current records for easy navigation.

 Data Viewing:

 Provides a detailed view of the selected database, including fields such as:

 Yearly Compensation, Gender, Age Group, Location, Job Title, Company Name

 The table view enables users to understand the structure and content of the dataset.

 Filtering Options:

 Users can filter records based on specific fields (e.g., compensation range, job title, or location).

 Filters allow for granular data exploration, enabling quick identification of trends and outliers.

 Data Manipulation Actions:

 View: Displays the full database for analysis and filtering.

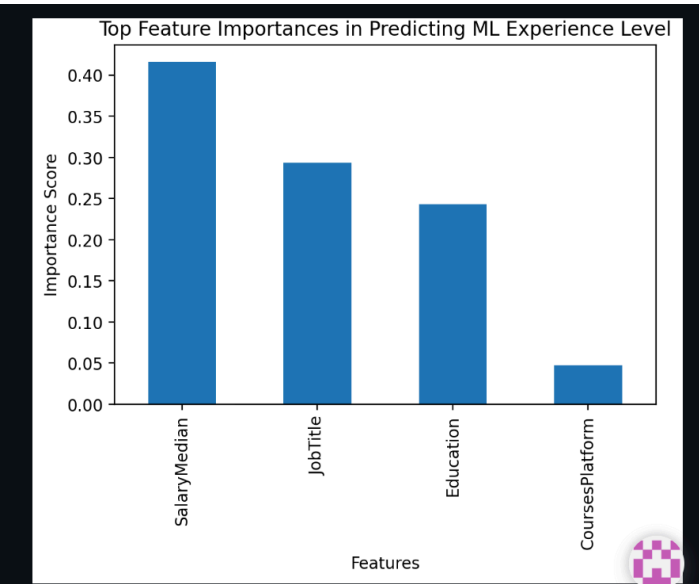 Add Entry: Users can input new data points directly into the database.
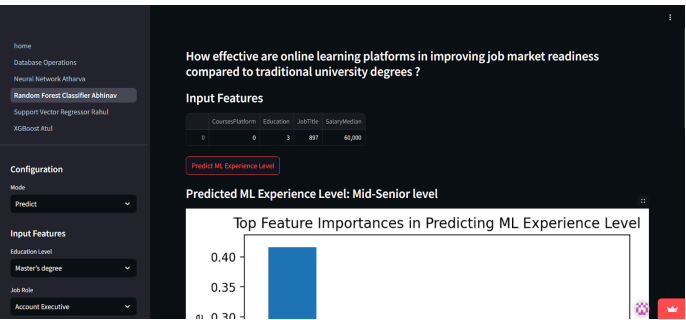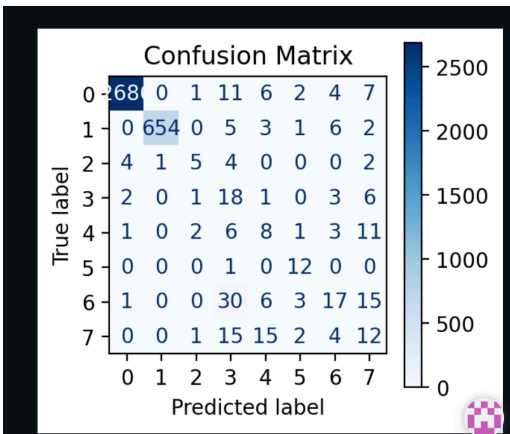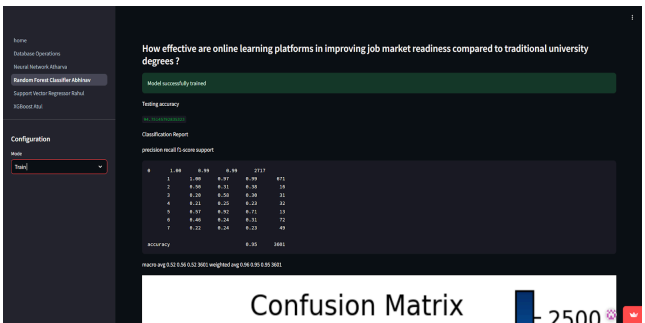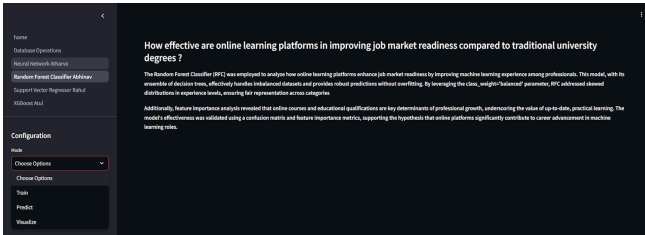
 Modify Entry: Existing records can be edited to update or correct information.

 Remove Entry: Unnecessary or outdated records can be deleted, ensuring the database remains relevant.

This page provides users with the capability to maintain and update datasets dynamically, ensuring that the predictions and insights generated by the application remain accurate and reflective of the latest data trends.
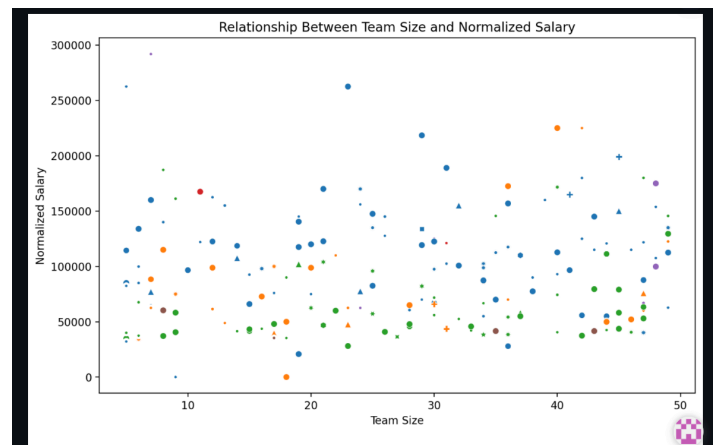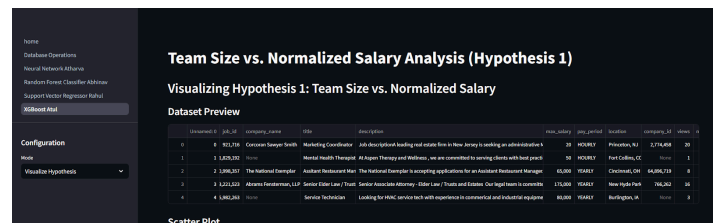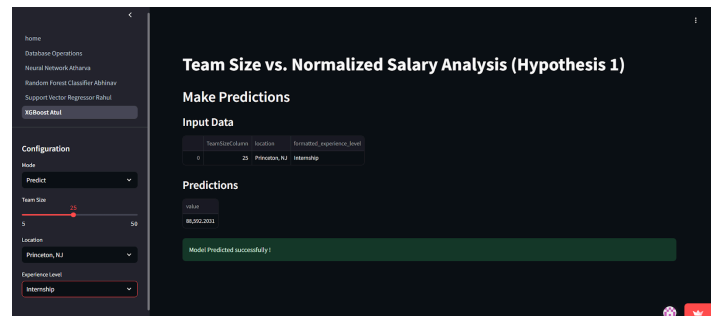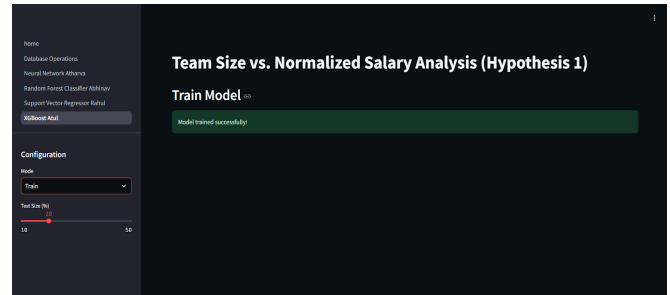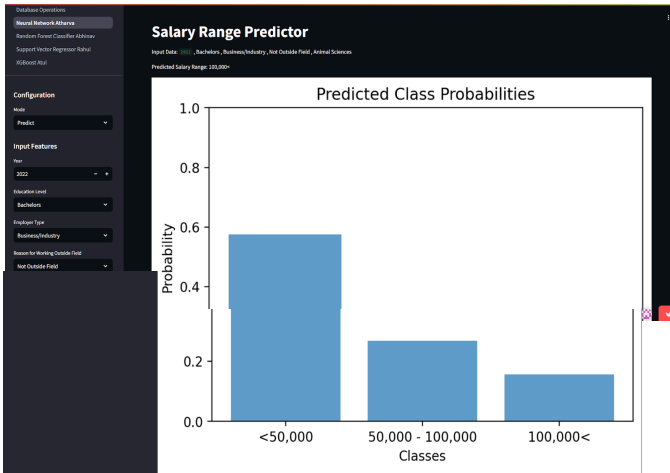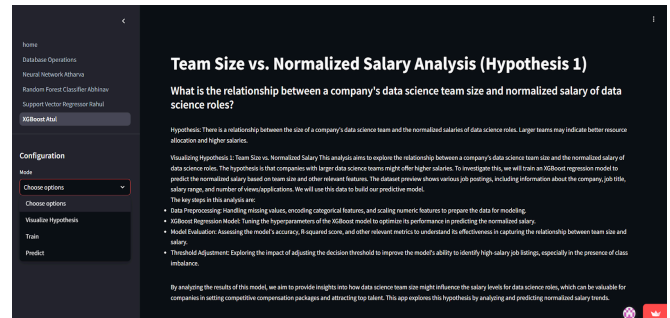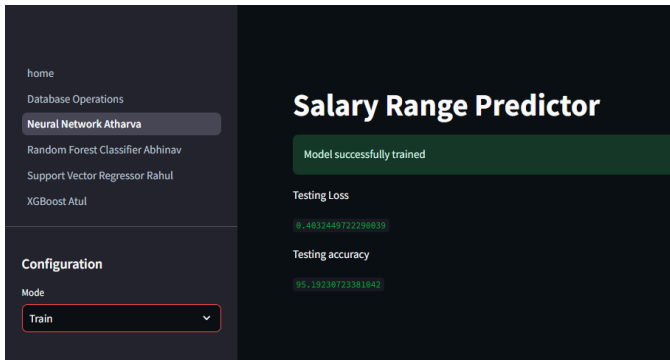
Abhinav Tembulkar:

This page, developed by Abhinav Tembulkar, addresses the question: "How effective are online learning platforms in improving job market readiness compared to traditional university degrees?" Utilizing a Random Forest Classifier, the model analyzes the impact of online courses on job market readiness by evaluating factors such as the course platform, education level, and job title. Users can interact with the page to train the model, make predictions, and visualize results. Key outputs include a classification report with accuracy, precision, recall, and F1 scores, a confusion matrix to evaluate prediction performance, and feature importance charts that reveal the significance of various factors in predicting job market readiness. The page provides a comprehensive analysis, enabling insights into the growing relevance of online learning platforms in career advancement.













Atharva Prabhu:

This page, developed by Atharva Prabhu, focuses on addressing the question: "Is it possible to predict expected salary range or employment status based on education major and demographic factors?" Using a Neural Network, this tool predicts the salary range based on user inputs such as educational major, employer type, education level, and demographic attributes like age and location. The salary range is categorized into three classes: below $50,000, between $50,000 and $100,000, and above $100,000. Users can interact with the page to either train the model on the provided dataset or predict salary outcomes based on their inputs. The results are visually displayed as a probability bar chart, showcasing the confidence levels for each salary range. This feature provides valuable insights into potential salary outcomes, enabling users to make informed decisions about their education and career paths.

**Atul Pandey:**

This page, developed by Atul Pandey, explores the hypothesis: "What is the relationship between a company's data science team size and the normalized salary of data science roles?" Using an XGBoost Regression model, the page predicts the normalized salary based on team size, job location, and experience level. The application provides three main functionalities: training the model, making predictions, and visualizing the results.

The training section allows users to preprocess data, optimize hyperparameters, and evaluate model performance using metrics such as R-squared scores. The prediction section accepts user inputs, such as team size and job-specific details, to estimate normalized salaries, which are displayed dynamically. The visualization section offers insights into the dataset, including a scatter plot showcasing the relationship between team size and salaries, enabling the identification of trends and patterns.

This page highlights how larger team sizes may indicate higher resource allocation, resulting in increased salary levels. It is a valuable tool for companies and professionals to analyze salary structures and make data-driven decisions for workforce planning and recruitment strategies.

Rahul Lotlikar:

This page, developed by Rahul Lotlikar, investigates the question: "Can we accurately predict an individual's job level within an organization based on demographic and professional characteristics such as age, experience level, and education?" Using Support Vector Regression (SVR), this analysis aims to predict job levels within organizations based on complex, high-dimensional data.

The training section allows users to preprocess the data, train the SVR model, and evaluate its performance using metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and $R^2$ Score. The prediction section enables users to input features like experience level, educational qualifications, and age to predict an individual's job level. The results include a numerical job level prediction, supplemented by a feature contribution plot that highlights the relative importance of each input feature in the prediction.

This page provides insights into the role of demographic and professional attributes in career progression, leveraging the SVR model's capability to capture non-linear relationships. It serves as a useful tool for organizations and individuals to analyze and optimize factors influencing career advancement within hierarchical structures.



REFERENCES

[1]A. Kon, "LinkedIn Job Postings," Kaggle Datasets, [Online]. Available: https://www.kaggle.com/datasets/arshkon/linkedin-job-postings.
[2] M. Exwell, "US Graduates Dataset," Kaggle Datasets, [Online]. Available: https://www.kaggle.com/datasets/mexwell/us-graduates.
[3] Kaggle, "Kaggle Survey 2019 Data," Kaggle Competitions, [Online]. Available: https://www.kaggle.com/competitions/kaggle-survey-2019/data.
[4]UOM190346A, "AI-Powered Job Market Insights," Kaggle Datasets, [Online].Available:https://www.kaggle.com/datasets/uom190346a/ai-powered-job-market-insights.
[5]U.S. Bureau of Labor Statistics, "Labor Force Statistics," [Online]. Available: https://www.bls.gov.[6]LinkedIn Corporation, "LinkedIn Data for Career Insights," LinkedIn Economic Graph Research, [Online]. Available: https://economicgraph.linkedin.com/.
[7]GitHub Repository, "Code Repository for ML Models and Application," [Online]. Available: https://github.com/APrabhu21/DIC_Project